

5. Slovensko-hrvatski paralelni korpus

Vesna Požgaj Hadži, Marko Tadić

Slovensko-hrvatski paralelni korpus bio je jedan je od projekata s područja humanističkih znanosti u sklopu Sporazuma o bilateralnoj slovensko-hrvatskoj suradnji u području znanosti i tehnologije između Ministarstva znanosti i tehnologije Republike Slovenije i Ministarstva znanosti i tehnologije Republike Hrvatske. »Težinu« projektu *Slovensko-hrvatski paralelni korpus* daje činjenica da je jedan od 28 prihvaćenih projekata (od 70 prijavljenih) koji je odobrio zajednički slovensko-hrvatski odbor krajem lipnja 1999. u trajanju od 2 godine (2000–2001). Riječ je zapravo o drugom dvojezičnom korpusu za slovenski jezik (prvi je *Slovensko-engleski korpus* Elan, v. Erjavec, 1999a, b). Rad na projektu započeo je početkom ak. god. 1999/2000.²² U poglavlju prikazujemo početak rada na projektu, njegove ciljeve i svrhu te upozoravamo na probleme vezane prije svega uz nedostatak digitalne grude za korpus.

5.1. Cilj i svrha projekta

Slovensko-hrvatski paralelni korpus bio je zamišljen kao ishodište za suvremena kontrastivna istraživanja dvaju genetski srodnih i susjednih jezika, koji su uspostavljanjem novih država doživjeli određene preinake između ostalog i zbog promjene državnog statusa. Ciljevi projekta su ovi:

- sastaviti usporedni korpus slovenskih i hrvatskih originala te odgovarajućih prijevoda (obostrani prijevodi),
- korpus sravniti (*align*) na razini rečeničnih prijevodnih ekvivalenta,
- omogućiti pristup korpusu na internetu putem web servisa.

Istraživanja korpusne lingvistike, koja posljednjih desetljeća 20. st. nadmašuju ostale lingvističke discipline, osobito u području leksikografije (jednojezičnih, višejezičnih i paralelnih korpusa, među koje je pripadao i ovaj projekt) čine temelj za čitav niz raznih vrsta lingvističkih proučavanja. Bilo je planirano da rezultati ovoga projekta nađu svoju primjenu u:

- kontrastivnim proučanjima (strukturnim i tipološkim) slovenskoga i hrvatskoga jezika na svim jezičnim razinama,
- leksikografskim proučanjima,

²² Kao partneri projekta pojavljuju se Filozofski fakultet Sveučilišta u Ljubljani i Filozofski fakultet Sveučilišta u Zagrebu. Projekt vode V. Požgaj Hadži i M. Tadić; u njemu s hrvatske strane sudjeluju I. Pranjković, V. Muhić-Dimanovski, B. Bekavac, S. Fulgos, K. Šojat; na slovenskoj strani na projektu rade V. Gorjanc, A. Skubic i Š. Vintar.

- Zbog zastarjelih slovensko-hrvatskih (srpskohrvatskih) rječnika (Juraničić³ 1986, ²1989)²³ projekt je trebao poslužiti kao temelj i poticaj za dvojezičnu leksikografiju. Slovensko-hrvatski rječnik i obratno danas postaje neophodnim normativnim priručnikom, čega su svjesni i mnogi nakladnici.
 - Slovensko-hrvatski paralelni korpus predstavlja i temelj za različita leksikografska i leksikološka proučavanja, npr. proučavanja »lažnih« prijatelja (rječnik slovensko-hrvatskih homonima), proučavanja terminologije (terminološki rječnici), proučavanja kolokacija itd.
- znanosti o prevođenju,
- Izrađeni korpus, objavljen na internetu, zamišljen je kao dragocjeno pomagalo studentima za vježbe iz prevođenja (slovenski-hrvatski i obratno) na Filozofskom fakultetu Sveučilišta u Ljubljani i Sveučilišta u Zagrebu, naročito zbog aktualnosti tekstova koje korpus sadrži. Naiime, korpus obuhvaća tekstove nastale posljednjih desetak godina 20. stoljeća, koji u oba jezika predstavljaju noviju jezičnu situaciju od one koja je postojala za vrijeme zajedničke države, a za razliku od rječnika nude i kontekst.
 - Izrađeni korpus dobro će doći i prevodiocima u obje države kao dopunski izvor informacija, posebice u posljednje vrijeme, kada se ponovno pojavljuju potrebe za kvalitetnim prijevodima, a kao što smo već rekli, dvojezični su nam rječnici zastarjeli i gotovo neupotrebljivi.
- razvoju jezičnih tehnologija,
- (automatsko) traženje termina i njihovih prijevodnih ekvivalenta,
 - strojno prevođenje i strojno potpomognuto prevođenje,
- didaktičkim/metodičkim proučavanjima,
- usustavljanju razlika između dvaju srodnih jezika koje su najčešće uzrokom pogrešaka u učenju/poučavanju jezika,
 - klasifikaciji interferencijskih pogrešaka u učenju/poučavanju hrvatskoga i slovenskoga jezika kao stranih i drugih, prije svega mislimo na fakultetsku razinu u obje zemlje te osnovnoškolsku u Sloveniji, u kojoj se od godine 2000. u 7. 8. i 9. razredu osnovne škole uvodi hrvatski kao izborni predmet (v. poglavlje 16),
 - izradi didaktičkih izvora (kontrastivnih udžbenika, priručnika, gramatička itd.) za sve razine i sve stupnjeve učenja obaju jezika.

²³ O nedostacima Juraničićevih rječnika v. Požgaj Hadži, 1998.

5.2. Sastavljanje korpusa

5.2.1. Korpusni parametri

Opseg korpusa na početku je definiran na milijun pojavnica (po 500.000 pojavnica za svaki jezik). Vremenski raspon korpusa obuhvaća tekstove nastale od 1990. do 2001, čime se pokušavalo približiti reprezentativnosti suvremene jezične situacije u oba jezika. Korpus ima ovu žanrovsку strukturu (v. tablicu 1): 40 % pripada stručnim tekstovima, 30 % publicističkim i po 15 % beletrističkim i znanstvenim tekstovima.

Tablica 1: *Struktura korpusa*

Tekstovi	Pojavnice
Beletristički tekstovi	15 %
Publicistički tekstovi	30 %
Stručni tekstovi	40 %
Znanstveni tekstovi	15 %

U početnim fazama prikupljanja tekstova za korpus preuzimali su se tekstovi isključivo u digitalnom zapisu. Ubrzo smo naišli na »tehničku prepreku« i uočili da takvih tekstova i takve žanrovske strukture nema dovoljno. Kako potpora projektu nije predviđala troškove za digitalizaciju tekstova, morali smo se ograničiti isključivo na tekstove koji su već postojali u digitalnome zapisu. Međutim, dostupnih tekstova u digitalnome zapisu – od kojih je jedan od paralelnih tekstova izvornik, a drugi njegov prijevod – jednostavno tada nije bilo dovoljno. To je dovelo do nemogućnosti ispunjenja planiranoga opsega korpusa i planirane žanrovske strukture. Također, ograničavajući čimbenik bio je i zahtjev da su tekstovi nastali nakon 1990. Tako postavljeni korpusni parametri pokazali su se prekonzervativnima i onemogućili su sastavljanje korpusa u željenome opsegu. Sakupljena grada do kraja 2000. godine obuhvaćala je međudržavne ugovore, beletristiku (romane i priče),²⁴ tehničku dokumentaciju i priručnike (promet, kulinarstvo, farmaceutika, elektrotehnika i sl.) te turističke brošure. Njezin je opseg oko 150.000 pojavnica u svakome jeziku.

24 Dramski i poetski tekstovi ne uzimaju se u korpus zbog svojih specifičnosti.

5.2.2. Obrada

Nad dijelom sakupljene građe probno je obavljena cijela obrada koja se sastojala od konverzije teksta u jedinstven XML zapis, priređivanja XML dokumenata za sravnjivanje i konačno samog postupka sravnjivanja. Konverzija teksta iz polaznih oblika zapisa (najčešće MS Word 97) obavljena je programom 2XML²⁵, koji je razvijen u Zavodu za lingvistiku Filozofskoga fakulteta Sveučilišta u Zagrebu u sklopu rada na *Hrvatskome nacionalnom korpusu*. Smjer je konverzije najčešće bio iz RTF u XML.

Nakon konverzije u XML obrada dokumenata sastojala se od segmentacije na rečenice <s> uz dodavanje ID oznaka rečenicama te provjere istovjetnosti broja odlomaka <p> kao pripreme za sravnjivanje.

5.2.3. Sravnjivanje

Za sravnjivanje je korištena DOS inačica programa Vanilla Aligner (Danielsson i Ridings 1997) koja se pokazala sasvim dostašnom za potrebe manjega paralelnog korpusa kao što je slovensko-hrvatski. Kako sam program i nema suviše fleksibilan ulazni zapis, XML dokumenti morali su mu se prilagoditi. U tom su procesu međutim zadržane XML oznake te su i one ušle u postupak sravnjivanja. Time je program dobio više redundantnih podataka i mogao je obaviti sravnjivanje na kvalitetniji način.

Probni uzorak čini osam međudržavnih ugovora između Republike Slovenije i Republike Hrvatske. Sami XML dokumenti prikazani su na slikama 1 i 2.

25 Sam program i njegovo funkcioniranje dijelom se prikazuje u Tadić (2000).

The screenshot shows a Microsoft Internet Explorer window displaying an XML document. The title bar reads "C:\HDPL2000\mt\veterin_si.XML - Microsoft Internet Explorer". The address bar shows the same path. The content area displays the XML code with some parts highlighted in blue, indicating they are links or selected text.

```
<BODY>
- <DIV0 type="main">
- <HEAD type="mainheading">
<S id="veterin_si.S1">SPORAZUM</S>
</HEAD>
- <HEAD type="subheading">
<S id="veterin_si.S2">MED VLADO REPUBLIKE
SLOVENIJE IN VLADO REPUBLIKE HRVAŠKE O
VETERINARSKEM SODELOVANJU</S>
</HEAD>
- <P>
<S id="veterin_si.S3">Vlada Republike Slovenije in
vlada Republike Hrvaške (v nadalnjem besedilu:
pogodbencii) sta se z željo, da bi pospešili promet z
živalmi in proizvodi živalskega izvora in z namenom,
da bi preprečili vnos kužnih bolezni živali in
proizvodov živalskega izvora, škodljivih za zdravje, in
da bi še naprej razvijali in učvrstili medsebojno
sodelovanje v veterinarstvu, ki temelji na načelih
enakosti, vzajemnega spoštovanja in skupne
blaginje,</S>
</P>
- <P>
<S id="veterin_si.S4">dogovorili o naslednjem:</S>
</P>
- <DIV1 type="sub">
- <HEAD type="article">
<S id="veterin_si.S5">1. člen</S>
</HEAD>
- <P>
<S id="veterin_si.S6">Uvoz in prevoz živali in
proizvodov živalskega izvora (v nadalnjem
besedilu: pošiljka) se opravlja, če so izpolnjeni
predpisani veterinarskosanitarni pogoji in če je
predhodno pridobljeno dovoljenje pristojnega
organa države uvoznice oziroma države, čez
ozemlje katere se prevaža pošiljka.</S>
</P>
- <P>
```

Slika 1: XML dokumenti u probnomre uzorku: slovenski ugovori

C:\HDPL2000\mt\veterin_hr.XML - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites History Mail

Address C:\HDPL2000\mt\veterin_hr.XML

```
- <BODY>
- <DIV type="main">
- <HEAD type="mainheading">
<S id="veterin_hr.S1">SPORAZUM</S>
</HEAD>
- <HEAD type="subheading">
<S id="veterin_hr.S2">IZMEĐU VLADE REPUBLIKE
SLOVENIJE I VLADE REPUBLIKE HRVATSKE O
VETERINARSKOJ SURADNJI</S>
</HEAD>
- <P>
<S id="veterin_hr.S3">Vlada Republike Slovenije i Vlada
Republike Hrvatske (u dalnjem tekstu: ugovorne
stranke), u želji da unaprijede promet životinja i
proizvoda životinjskog podrijetla i s namjerom da
zapriječe unošenje zaraznih bolesti životinja i po
zdravlje štetnih proizvoda životinjskog podrijetla, kao
i da dalje razvijaju i učvrste međusobnu suradnju u
veterinarstvu utemeljenu na načelima jednakosti,
uzajamnog štovanja i zajedničke dobrobiti,</S>
</P>
- <P>
<S id="veterin_hr.S4">dogovorile su se kako
slijedi:</S>
</P>
- <DIV1 type="sub">
- <HEAD type="article">
<S id="veterin_hr.S5">Članak 1.</S>
</HEAD>
- <P>
<S id="veterin_hr.S6">Uvoz i provoz životinja i
proizvoda životinjskog podrijetla (u dalnjem
tekstu: pošiljka) obavlja se ako je udovoljeno
propisanim veterinarsko-sanitarnim uvjetima i ako
je prethodno pribavljeno odobrenje nadležnog
tijela zemlje uvoznice odnosno zemlje preko čijeg
se teritorija pošiljka provozi.</S>
</P>
- <P>
```

Done My Computer

Slika 2: XML dokumenti u probnom uzorku: hrvatski ugovori

5.3. Rezultati probnoga uzorka

Preliminarna statistika probnoga uzorka osam međudržavnih ugovora između Republike Slovenije i Republike Hrvatske po XML elementima prikazana je u tablici 4:

Tablica 2: Broj *<p>*, *<s>* i *<w>* elemenata u probnom uzorku

	Hrvatski	Slovenski
Odlomaka <i><p></i>	522	522
Rečenica <i><s></i>	891	891
Pojavnica <i><w></i>	13.549	13.307

Valja uočiti identičnost broja odlomaka i rečenica koja se pojavljuje zbog naravi pravnih tekstova. Tekstovi na oba jezika objavljeni su kao originali i uskladeni već u izvorniku. Stoga su sva strojna sravnjivanja u tom probnom uzorku oblika 1–1. Gotovo identično ponašanje može se naći u drugim paralelnim korpusima pravnih tekstova (v. Gamper, 2000). Dapače, kad se god pri pregledu sravnjivanja naišlo na 2–1 ili 1–2 sravnjivanje, to je redovito bio znak pogreške u ulaznim podacima tj. pogreške u modulu za segmentaciju rečenica.

Rezultat sranjivanja Vanillom prikazan je na slici 3, a u sažetijem zapisu na slici 4.

```

veterni_all.txt - Notepad
File Edit Search Help

*** Lämk: 1 - 1 ***
<HEAD type="main"> <HEAD type="mainheading"> <S id="veterin_hr_S1"> IZME#272;U ULADE REPUBLIKE SLOVENIJE I ULADE REPUBLIKE
HRVATSKE O VETERINARSKOJ </S> </HEAD> .EOS
<HEAD type="subheading"> <S id="veterin_si_S2"> MED ULADO REPUBLIKE SLOVENIJE IN ULADO REPUBLIKE
HRUR#352;KE O VETERINARSKEM SODELOVANJU </S> </HEAD> .EOS
.EOP

*** Lämk: 1 - 1 ***
<P> <S id="veterin_hr_S3"> Ulada Republike Slovenije i Ulada Republike Hrvatske (u daljinjem tekstu:
ugovorne stranke), u #382;elji da unaprjede promet #382;ivotinja i proizvoda #382;ivotinskog
podrjetja i s namjerom da zaprijet#269;e uno#33;enje zaravnih bolesti #382;ivotinja i po zdravlj
&#353;tetnih proizvoda #382;ivotinskog podrijetla, kao i da dale razvijaju i ut#269;urste
me&#273;usobnu suradnju u veterinarstvu utemeljenu na na#269;elima jednakosti, uzajamnog &#353;tovanja
i zajednic#269;e dobrobiti, </S> </P> .EOS
<P> <S id="veterin_si_S3"> Ulada Republike Slovenije in Ulada Republike Hrvat#353;ke (u nadalnjem
besedilu: pogodbencii) sta se z #382;eljo, da bi pospet#353;ili promet z #382;idalni in prizvodi
&#382;ivalistega izvora in z nanenom, da bi preprek#269;ili unos kuka#382;ih bolezni #382;ulovi in
proizvodou #382;ivalskoga izvora, #353;Kodljivih za zdravje, in da bi #382;naprej razviali in
u&#269;urstili medsebojno sodelovanje u veterinarstvu, ki temelji na na#269;elih enakosti, uzajemnega
spotet#353;tovanja in skupne blaginje, </S> </P> .EOS
.EOP

*** Lämk: 1 - 1 ***
<P> <S id="veterin_hr_S4"> dogovorile su se kako slijedi: </S> </P> .EOS
<P> <S id="veterin_si_S4"> dogovorili o naslednjem: </S> </P> .EOS
.EOP

*** Lämk: 1 - 1 ***
<DI1 type="sub"> <HEAD type="article"> <S id="veterin_hr_SS"> &#268;lanak 1. </S> </HEAD> .EOS
<DI1 type="sub"> <HEAD type="article"> <S id="veterin_si_SS"> 1. &#268;len </S> </HEAD> .EOS

```

Slika 3: Rezultat stvarnjivanja Vanillom

The screenshot shows a Microsoft Internet Explorer window displaying XML code. The title bar reads "C:\temp7\veterin_al2.xml - Microsoft Internet Explorer". The address bar shows the file path "C:\temp7\veterin_al2.xml". The content area contains the following XML code:

```
- <body>
<link xtargs="veterin_hr.S1 veteran_si.S1" />
<link xtargs="veterin_hr.S2 veteran_si.S2" />
<link xtargs="veterin_hr.S3 veteran_si.S3" />
<link xtargs="veterin_hr.S4 veteran_si.S4" />
<link xtargs="veterin_hr.S5 veteran_si.S5" />
<link xtargs="veterin_hr.S6 veteran_si.S6" />
<link xtargs="veterin_hr.S7 veteran_si.S7" />
<link xtargs="veterin_hr.S8 veteran_si.S8" />
<link xtargs="veterin_hr.S9 veteran_si.S9" />
<link xtargs="veterin_hr.S10 veteran_si.S10" />
<link xtargs="veterin_hr.S11 veteran_si.S11" />
<link xtargs="veterin_hr.S12 veteran_si.S12" />
<link xtargs="veterin_hr.S13 veteran_si.S13" />
<link xtargs="veterin_hr.S14 veteran_si.S14" />
<link xtargs="veterin_hr.S15 veteran_si.S15" />
<link xtargs="veterin_hr.S16 veteran_si.S16" />
<link xtargs="veterin_hr.S17 veteran_si.S17" />
<link xtargs="veterin_hr.S18 veteran_si.S18" />
<link xtargs="veterin_hr.S19 veteran_si.S19" />
<link xtargs="veterin_hr.S20 veteran_si.S20" />
<link xtargs="veterin_hr.S21 veteran_si.S21" />
<link xtargs="veterin_hr.S22 veteran_si.S22" />
<link xtargs="veterin_hr.S23 veteran_si.S23" />
<link xtargs="veterin_hr.S24 veteran_si.S24" />
<link xtargs="veterin_hr.S25 veteran_si.S25" />
<link xtargs="veterin_hr.S26 veteran_si.S26" />
<link xtargs="veterin_hr.S27 veteran_si.S27" />
<link xtargs="veterin_hr.S28 veteran_si.S28" />
<link xtargs="veterin_hr.S29 veteran_si.S29" />
<link xtargs="veterin_hr.S30 veteran_si.S30" />
<link xtargs="veterin_hr.S31 veteran_si.S31" />
<link xtargs="veterin_hr.S32 veteran_si.S32" />
<link xtargs="veterin_hr.S33 veteran_si.S33" />
<link xtargs="veterin_hr.S34 veteran_si.S34" />
<link xtargs="veterin_hr.S35 veteran_si.S35" />
<link xtargs="veterin_hr.S36 veteran_si.S36" />
<link xtargs="veterin_hr.S37 veteran_si.S37" />
<link xtargs="veterin_hr.S38 veteran_si.S38" />
</body>
```

Slika 4: Rezultat sravnjivanja Vanillom u sažetijem zapisu

U poglavlju je prikazan početak rada na projektu *Slovensko-hrvatski paralelni korpus* koji je odobrilo i finansiralo i slovensko i hrvatsko Ministarstvo znanosti i tehnologije te je time omogućena suradnja stručnjaka dvaju istovrsnih fakulteta. U prvoj godini rada na projektu (2000) sakupljeni tekstovi, različitih funkcionalnih stilova i različitih područja, zapravo predstavljaju kompromis između očekivane uporabnosti korpusa s jedne strane i same dostupnosti tekstova u digitalnom zapisu s druge strane. Naime, prvih desetak godina nakon osamostaljenja Republike Slovenije i Republike Hrvatske, bilo je uočljivo nepostojanje prijevoda sa slovenskoga na hrvatski i hrvatskoga na slovenski jezik (osobito u digitalnom obliku), što je rezultiralo značajnim teškoćama oko prikupljanja građe za korpus. Kad bude zaokružen u planiranome opsegu i kad postane dostupan putem interneta, *Slovensko-hrvatski paralelni korpus* omogućiće niz različitih jezikoslovnih istraživanja: kontrastivnih, leksikografskih, didaktičko/metodičkih, istraživanja vezanih uz znanost o prevodenju te razvitak jezičnih tehnologija za oba jezika. Potencijalni su korisnici korpusa, osim istraživača i studenata slovenistike i kroatistike, također i prevoditelji kojima će korpus biti dostupan putem www-a kao dopunski izvor informacija, osobito u situaciji kada dvojezičnih hrvatsko-slovenskih (i obratno) rječnika zapravo nema.²⁶

26 Korpus, nažalost, zbog teškoća koje smo predviđali već u prvoj godini rada na projektu (nedostatak tekstova u digitalnom obliku) nije završen i tako nedostupan korisnicima.